



Portable Implementation of the p2z Benchmark using Alpaka

Cong Wang

BACKGROUND

- Today's computing in HEP:
 - code written for the x86 platform is pretty much guaranteed to run everywhere in the world, from computing centers using batch systems to our own laptops.
 - GPUs and other accelerators provide more processing power for the same energy consumption as with x86-based supercomputers. It can lead to designing algorithms and implementing them for specific hardware platforms and combinations, and making it very difficult to use not only one but several of these platforms.
- Motivation:
 - investigate solutions for portability techniques that will allow the coding of an algorithm once, and the ability to execute it on a variety of hardware products from many vendors, especially including accelerators.

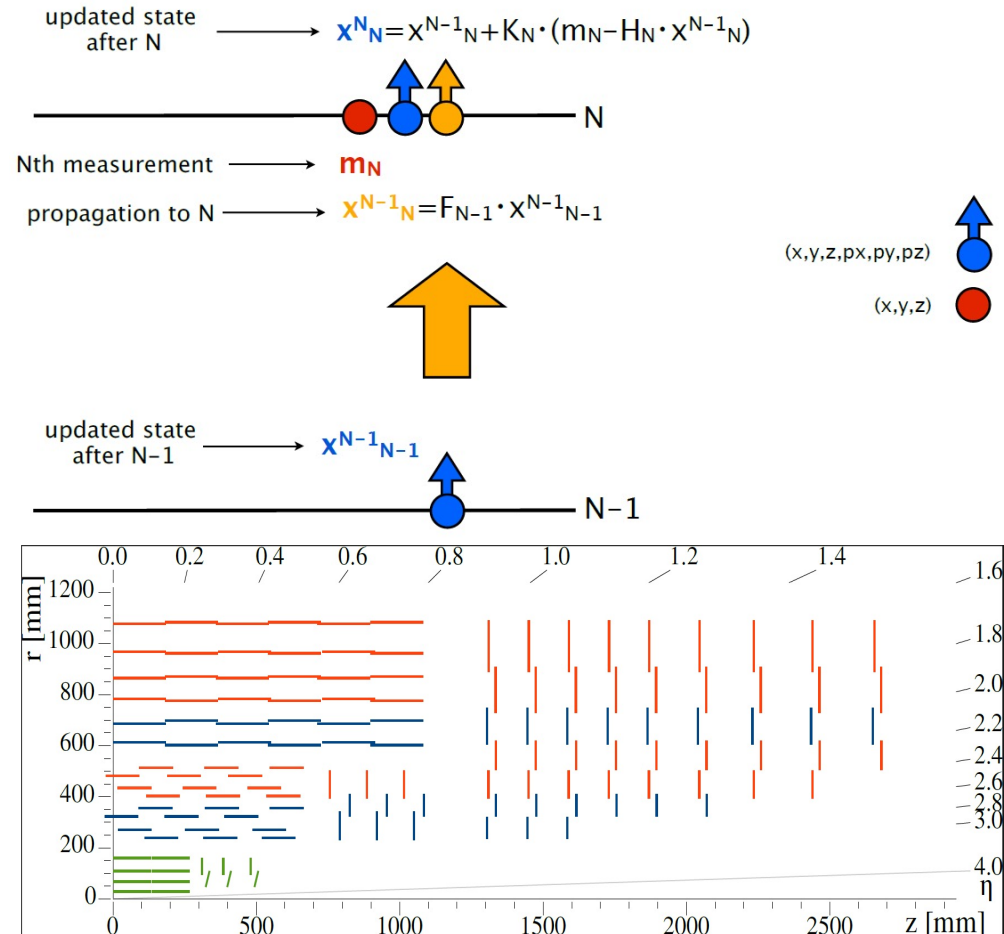
BACKGROUND

- Propagate to z (p2z):

- the layers at a fixed z positions (endcap disks). The particles are propagated in a homogeneous magnetic field with direction parallel to the z axis.
- each layer is propagated from layer N-1 to N and updated the track parameters based on the hit (measurement) located on layer N (mN) using the Kalman Filter.
- collision events inside a particle detector are independent of each other.
- the various tracks created within an event can be built and fitted independently.
- tracks are grouped together to propagate in lockstep from one layer of the detector to the next.

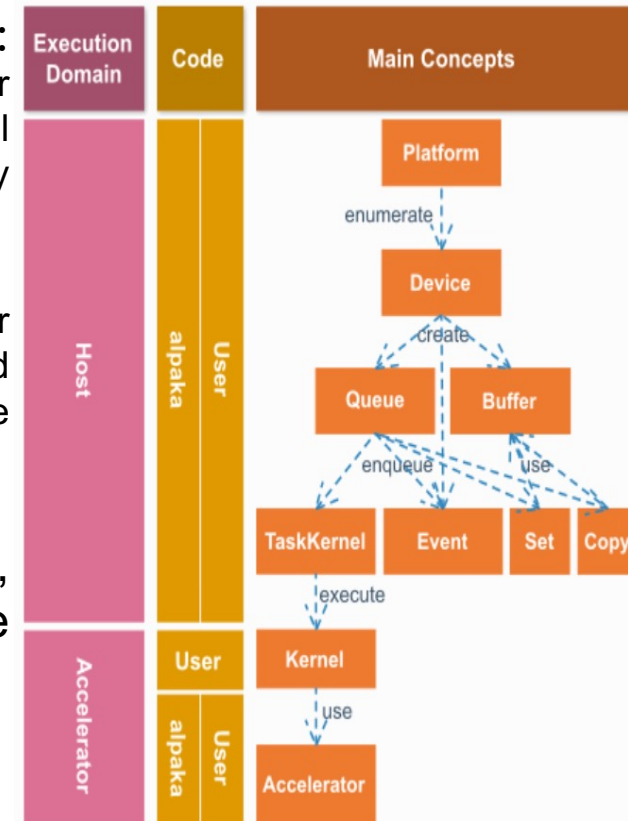
- Parameters:

- batch size (bsize), number of layers (nlayer), events (nevts), tracks (ntrks), and batches (nb=ntrks/bsize)



ALPAKA

- **Abstraction Library for Parallel Kernel Acceleration:**
 - defines and implements an abstract interface for the *hierarchical redundant parallelism* model. The model exploits task- and data-parallelism as well as memory hierarchies at all levels of current multi-core architectures.
 - provides back-ends for *CUDA*, *OpenMP*, *TBB* and other methods. The policy-based C++ template interface provided allows for straightforward user-defined extension of the library to support other accelerators.
- Sustainable, heterogeneous, maintainable, testable, optimizable, extensible, data structure agnostic



IMPLEMENTATION

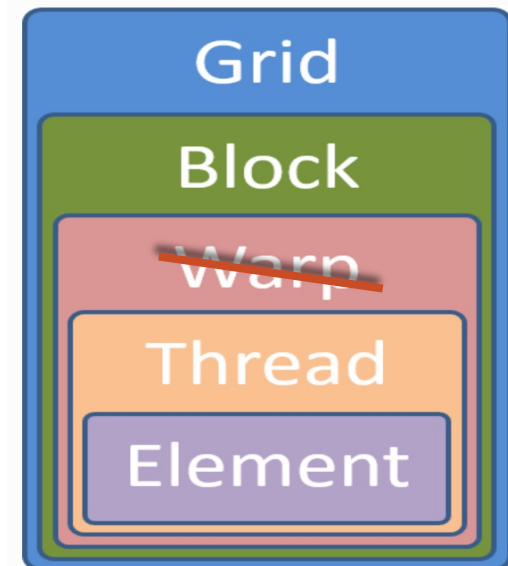
- **Propagate-toz-test_alpaka_CPU:**

- a universal grid-block and block-thread definition for **CpuThreads**, **Omp2threads**, **Omp2Blocks**, and **TbbBlocks** as portable accelerators.
- each event & batch of tracks is distributed to one thread or block.
- SIMD inside computing functions. **Each thread processes # of batch size elements** which could be treated at once due to the same instructions and no data dependency.
- based on the accelerator, explicitly define the number of resources, loop over till the last event & batch of tracks.


threadIdx.x, threadIdx.y
= defined parallel threads

blockIdx.x * blockIdx.y
= defined parallel blocks

elementsperthread = batch size

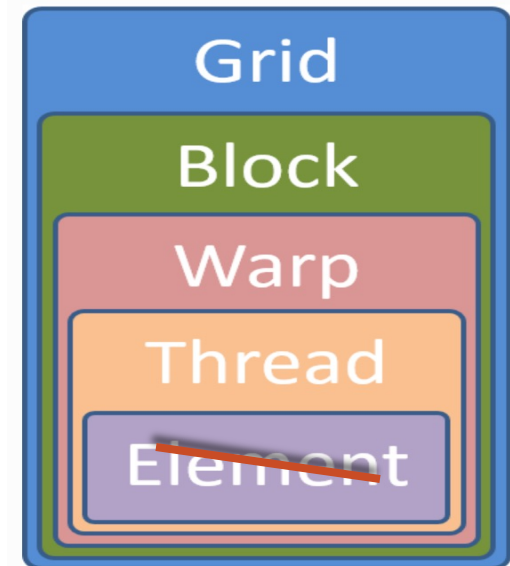


IMPLEMENTATION

- **Propagate-toz-test_alpaka_GPU:**
 - same universal grid-block and block-thread definition as CPU for **CudaRt** accelerator. 
 - each event & batch of tracks is distributed to one block, one thread are responsible for one track.
 - based on the accelerator, explicitly define the number of resources, loop over till the last event & batch of tracks.
 - one to multiple streams, asynchronous memory transfer.

$\text{threadIdx.x} * \text{threadIdx.y}$
= defined parallel threads

$\text{blockIdx.x} * \text{blockIdx.y}$
= defined parallel blocks



CODE EXAMPLE

- Kernel:

```

610 struct alpakaKernel
611 {
612 public:
613     template<typename TAcc>
614     ALPAKA_FN_ACC auto operator()(
615         TAcc const& acc,
616         MPTRK* trk,
617         MPHIT* hit,
618         MPTRK* outtrk,
619         const int stream
620     ) const -> void
621     {
622         using Dim = alpaka::Dim<TAcc>;
623         using Idx = alpaka::Idx<TAcc>;
624         using Vec = alpaka::Vec<Dim, Idx>;
625
626         Vec const threadIdx = alpaka::getIdx<alpaka::Block, alpaka::Threads>(acc);
627         Vec const threadExtent = alpaka::getWorkDiv<alpaka::Block, alpaka::Threads>(acc);
628         Vec const blockIdx = alpaka::getIdx<alpaka::Grid, alpaka::Blocks>(acc);
629         Vec const blockExtent = alpaka::getWorkDiv<alpaka::Grid, alpaka::Blocks>(acc);
630
631         auto & errorProp = alpaka::declareSharedVar<MP6x6F, __COUNTER__>(acc);
632         auto & temp = alpaka::declareSharedVar<MP6x6F, __COUNTER__>(acc);
633         auto & inverse_temp = alpaka::declareSharedVar<MP3x3, __COUNTER__>(acc);
634         auto & kGain = alpaka::declareSharedVar<MP3x6, __COUNTER__>(acc);
635         auto & newErr = alpaka::declareSharedVar<MP6x6SF, __COUNTER__>(acc);
636
637         int ie_range;
638         if(stream == num_streams){ ie_range = (int)(nevts%num_streams);}
639         else{ie_range = (int)(nevts/num_streams);}
640
641         for (size_t ie=blockIdx[0]; ie<ie_range; ie+=blockExtent[0]){ //loop for TbbBlocks & Omp2Blocks & GPU
642             for (size_t ib=blockIdx[1]; ib<nb; ib+=blockExtent[1]){ //loop for TbbBlocks & Omp2Blocks & GPU
643                 const MPTRK* btracks = bTk(trk, ie, ib);
644                 MPTRK* obtracks = bTk(outtrk, ie, ib);
645                 for( size_t layer=0; layer<nlayer;++layer){
646                     const MPHIT* bhits = bHit(hit, ie, ib, layer);
647                     //struct MP6x6F errorProp, temp;
648                     propagateToZ(&btracks).cov, &btracks).par, &btracks).q, &bhits).pos, &btracks).cov, &btracks).par, &errorProp, &temp, acc); // vectorized function
649                     KalmanUpdate(&btracks).cov, &btracks).par, &bhits).cov, &bhits).pos, &inverse_temp, &kGain, &(newErr), acc);
650                 }
651             }
652         }
653     }
654 };

```

For GPU & CPU

Events →

Bunch of Tracks ↓

Block (0,0) Block (1,0) ...

Block (0,1) Block (1,1)

...

CODE EXAMPLE

- Device (accelerator) functions:

```

557 template< typename TAcc>
558 inline void ALPAKA_FN_ACC propagateToZ(const MP6x6SF* inErr, const MP6F* inPar, const MP1I* inChg, const MP3F* msP, MP6x6SF* outErr, MP6F* outPar, MP6x6F* errorProp, MP6x6F* temp, TAcc const & acc) {
559     using Dim = alpaka::Dim<TAcc>;
560     using Idx = alpaka::Idx<TAcc>;
561     using Vec = alpaka::Vec<Dim, Idx>;
562
563     Vec const ElementExtent = alpaka::getWorkDiv<alpaka::Thread, alpaka::Elems>(acc);
564     Vec const threadIdx = alpaka::getIdx<alpaka::Block, alpaka::Threads>(acc);
565     Vec const threadExtent = alpaka::getWorkDiv<alpaka::Block, alpaka::Threads>(acc);
566
567     for(size_t i=threadIdx[0]; i<bsize; i+=threadExtent[0]){
568         #pragma omp simd
569         for(size_t ele=0; ele<ElementExtent[0]; ele++){
570             size_t it = ele + ElementExtent[0] * i;
571             const float zout = z(msP, it);
572             const float k = q(inChg, it)*kfact;//100/3.8;
573             const float deltaZ = zout - z(inPar, it);
574             const float pt = 1./ipt(inPar, it);
575             const float cosP = cosf(phi(inPar, it));
576             const float sinP = sinf(phi(inPar, it));
577             const float cosT = cosf(theta(inPar, it));
578             const float sinT = sinf(theta(inPar, it));

```

Inside Block (0,0)

For GPU

For CPU

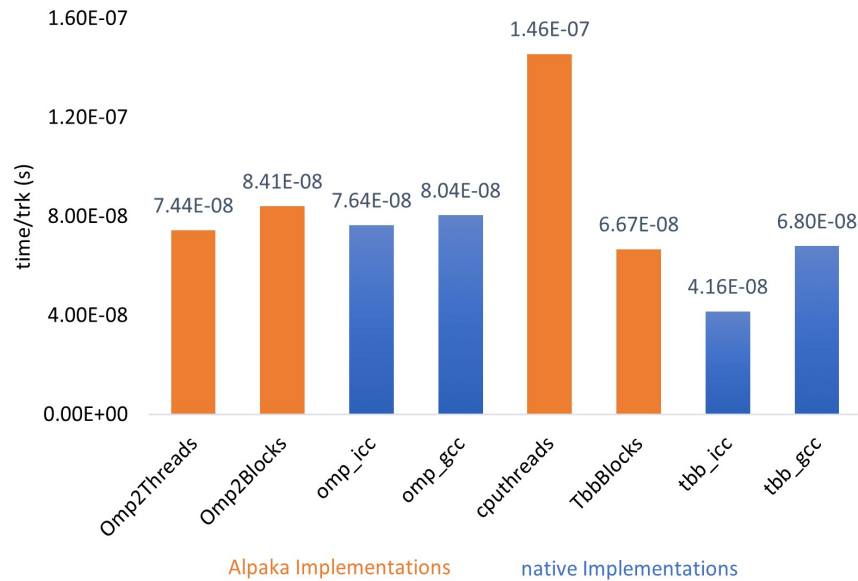
The diagram illustrates the execution model for GPU and CPU. For GPU, multiple threads (e.g., Thread (0,0), Thread (0,1)) are shown, each with a vertical arrow labeled 'Tracks' pointing downwards, indicating parallel execution. For CPU, a single thread (Thread (0,0)) is shown, containing multiple elements (e0, e1, ..., eN) with a vertical arrow labeled 'Tracks' pointing downwards, indicating sequential execution of elements within a thread.

COMPILATION

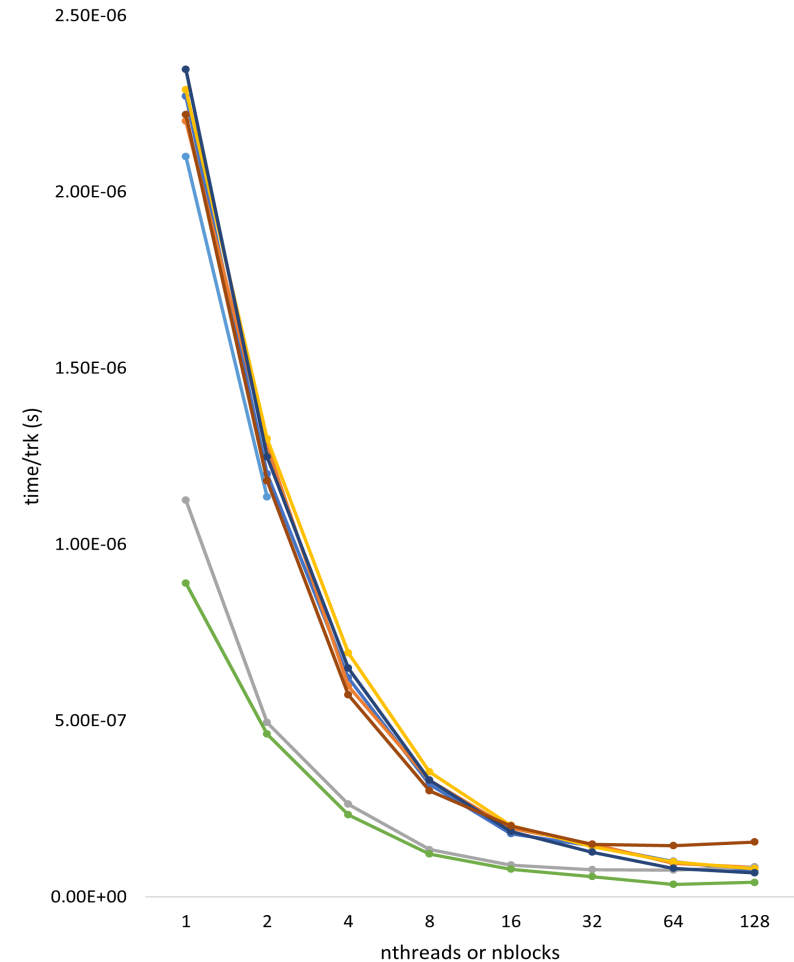
- Parameters:
 - bsize: 32, ntrks: 9600, nevts: 100, nb: 300, smear: 0.1, nlayer: 20, NITER: 5, Stream: 5
- Computing Architecture:
 - Intel Xeon Gold 6148@2.40 GHz on Apollo for CPU, 80 cores in total.
 - V100 16 GB GPU on Wilson Cluster for GPU implementation testing.
- Compilers:
 - GCC, NVCC
- Cmake option:
 - DCMMAKE_CXX_COMPILER=g++
DCMAKE_C_COMPILER=gcc
DCMAKE_CUDA_COMPILER=nvcc
 - DALPAKA_ACC_CPU_B_SEQ_T_SEQ_ENABLED -
DALPAKA_ACC_CPU_B_OMP2_T_SEQ_ENABLED -
DALPAKA_ACC_CPU_B_SEQ_T_OMP2_ENABLED -
DALPAKA_ACC_CPU_B_SEQ_T_THREADS_ENABLED -
DALPAKA_ACC_CPU_B_TBB_T_SEQ_ENABLED -
DALPAKA_ACC_GPU_CUDA_ENABLE -

RESULTS CPU

Peak Performance

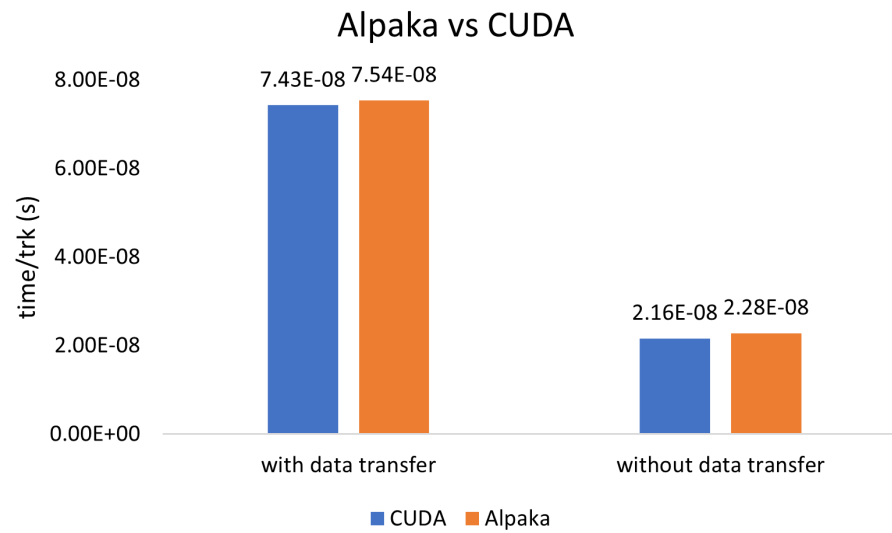


Scalability & Benchmarking



Omp2Threads Omp2Blocks omp_icc omp_gcc
 TbbBlocks tbb_icc tbb_gcc cputhreads

RESULTS GPU



CONCLUSION

- Alpaka Experience and Portability Matrix

Learning	Code Conversion	Building level	Hardware map	Feature
Easy to learn Good docs Lack of examples	Convertible with little more efforts	No major changes CMake provided	CPU, Nvidia GPU, AMD GPU No other supported	Reduction, Atomic Kernel Concurrency
Debugging	User Support	Sustainability	Interoperability	Performance
Easy to debug	Discussing thread Small community	Life cycle unpredictable	Mix with compiler directives and CUDA API	Minor loss or equivalent

REFERENCE

- E. Zenker, R. Widera, G. Juckeland et al., *Porting the Plasma Simulation PIConGPU to Heterogeneous Architectures with Alpaka*, [video link \(39 min\)](#), [slides \(PDF\)](#), [DOI:10.5281/zenodo.6336086](#)
- Alpaka: <https://alpaka.readthedocs.io/en/latest/index.html>
- Lantz, Steven, et al. "Speeding up particle track reconstruction using a parallel Kalman filter algorithm." *Journal of Instrumentation* 15.09 (2020): P09030.
- Bhattacharya, Meghna, et al. "Portability: A Necessary Approach for Future Scientific Software." *arXiv preprint arXiv:2203.09945* (2022).

QUESTIONS

